As per the Latest Syllabus of Anna University, Chennai (Regulation-2017)

# BIG DATA ANALYTICS

## For B.E./B.Tech VI Semester IT & VII Sem CSE Branches (Professional Elective)



**ARS PUBLICATIONS**
CHENNAI.

**S. ARUNPRASATH**
**K. SRIRAM KUMAR**
**P. KRISHNA SANKAR**

# PREFACE

This book "Big Data Analytics" is to know about the fundamental concepts of big data, streams and analytics, with various tools and practices in real world. It contributes an impression towards big data programming concepts of R and Python. It provides a preliminary study to access and perform analytics on huge volume of data. It affords procedural footsteps and study over NoSQL, Twitter data analytics and Wikipedia blog.

**Unit I:** Introduction towards evolution, best practices and characteristics of Big Data. Outline about use cases on Bigdata storage and architecture, real world Hadoop analytics mechanism available currently.

**Unit II:** Outline towards clustering, K-means and procedural steps in cluster construction. Classification and its core mechanism decision tree, Naïve Bayes are systematically briefed with R program.

**Unit III:** Transient awareness on association rules, Apriori algorithm and recommendation system. Brief knowledge over detecting candidate rules and collaborative, Content based, Knowledge based and Hybrid Recommendation with applications.

**Unit IV:** Contributes a knowledge on trendy Stream computing and its architecture. Real world analytics like Sentiment analysis on Twitter, Stock market prediction and Graph analytics are briefed with understanding them over stream computing.

**Unit V:** Provides a study over NoSQL and various real-world methodology. Various case studies on Hive and Hadoop architecture used in Twitter, E-Commerce and Blogs are briefed. It provides introduction towards R programming and its available function to perform data analytics.

# Contents

**UNIT III**

**ASSOCIATION AND RECOMMENDATION SYSTEM**

**Unit IV**

**Stream Memory**